# PowerTCP

## Pushing the Performance Limits of Datacenter Networks

Vamsi Addanki, Oliver Michel, Stefan Schmid
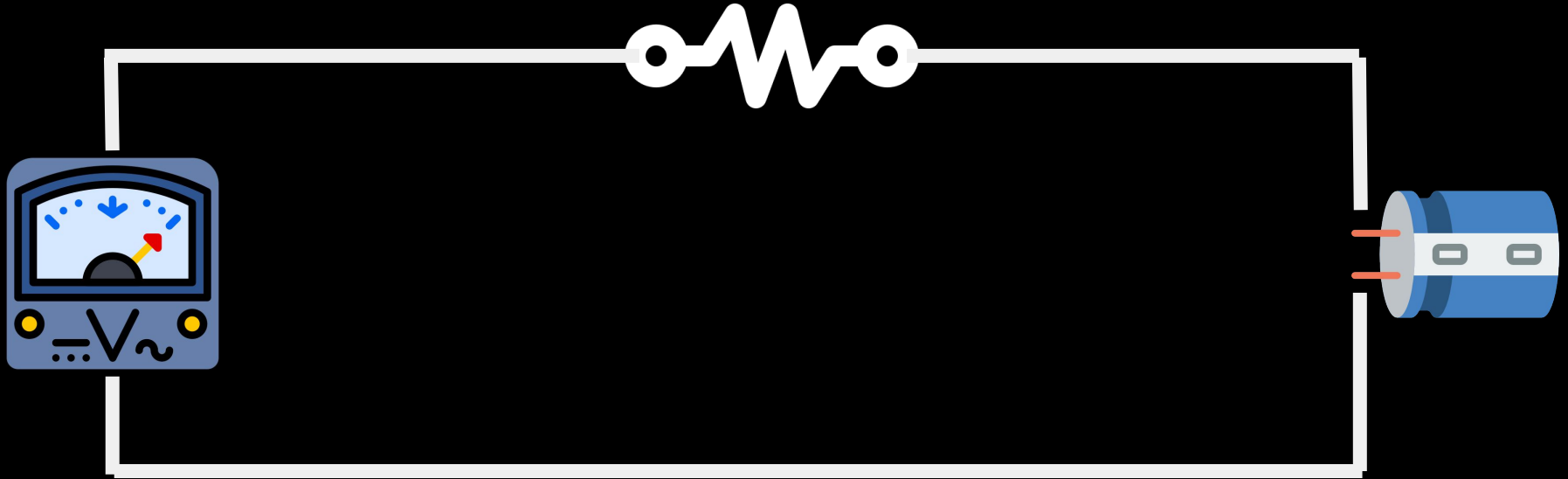
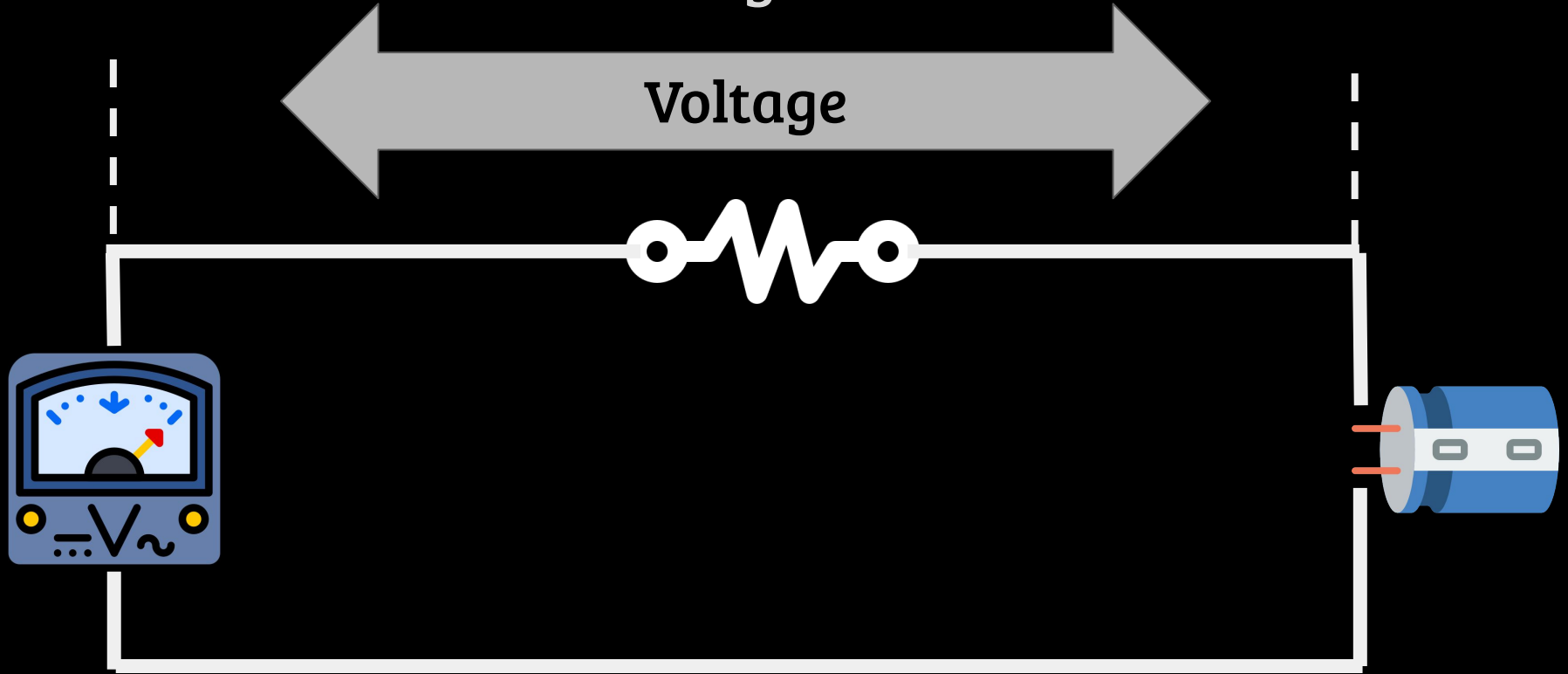# Brief context of electrical systems

# Brief context of electrical systems

**Voltage**

PowerTCP

# Brief context of electrical systems

**Voltage**

**Current**

Brief context of electrical systems

Voltage

Current

**Power = Voltage x Current**

# Analogy to networked systems

**Bottleneck Queue**

**Sender**

**Receiver**

**Feedback**

# Analogy to networked systems

**Voltage = BDP + queue length**

# Analogy to networked systems

Voltage = BDP + queue length

Current = Rate

Analogy to networked systems

Voltage = BDP + queue length

Current = Rate

Power?

Upnext... stay tuned!

PowerTCP

9

# PowerTCP in a Nutshell

- **Power**-based congestion control
- Quickly reacts to congestion **without losing throughput**
- Rapidly converges **within 1 RTT**
- Fair and **asymptotically** stable
- Reduces FCTs for short flows **by up to 90%**

POWERTCP

# How do we measure Power?

**The debate over congestion signals**

Microsoft says ECN is better [dctcp]

Google says delay is simple and effective [Timely, Swift]
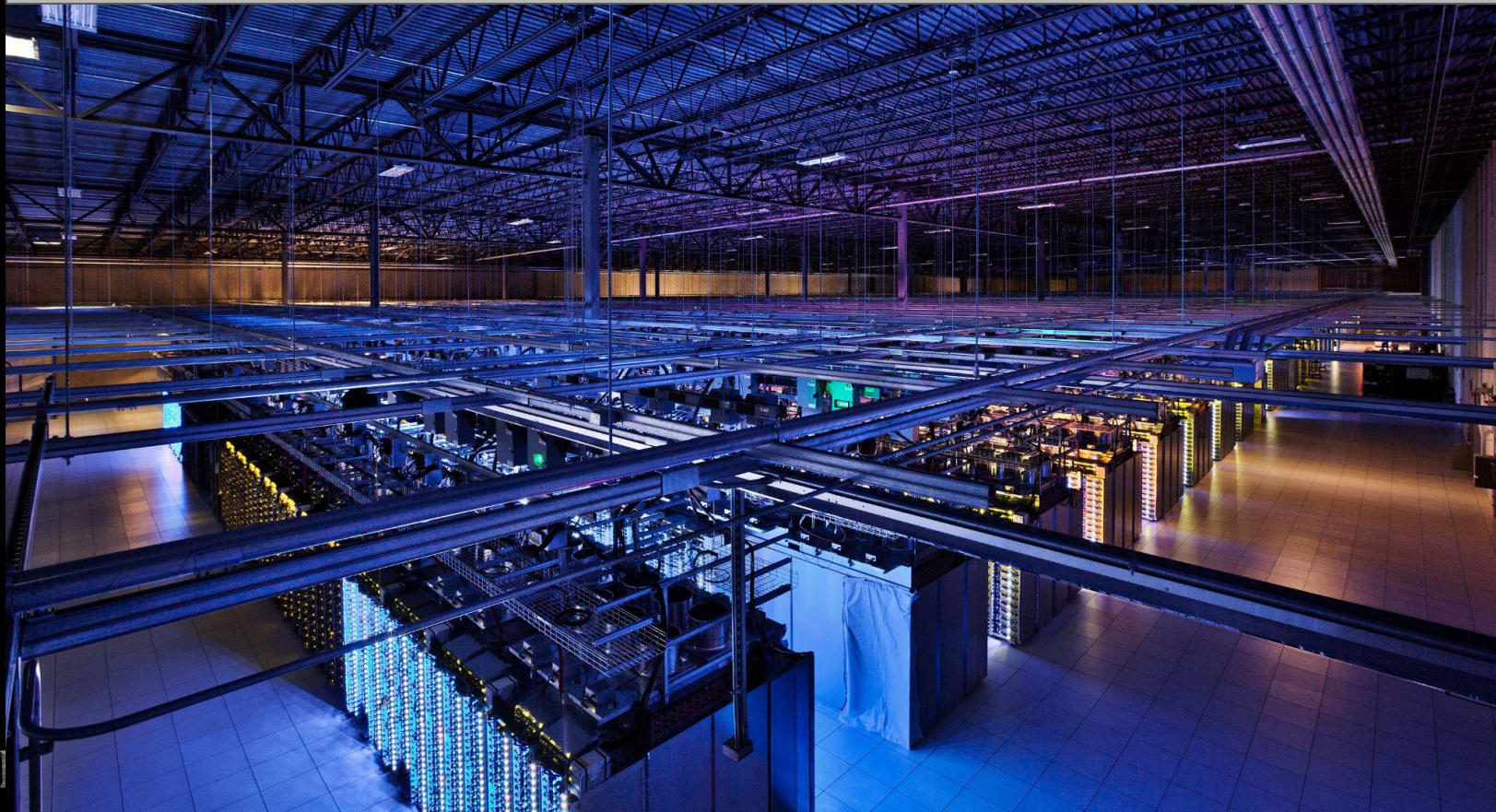
Alibaba says INT is accurate [HPCC]

ECN, Delay or INT are essential

What matters more: what we do with it

POWERTCP

~~The debate over feedback signals~~

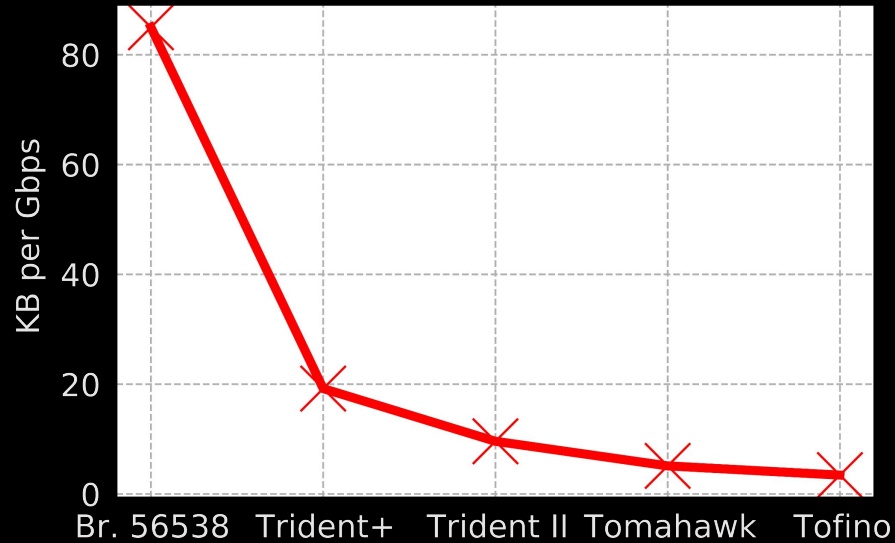A debate over how to use the feedback

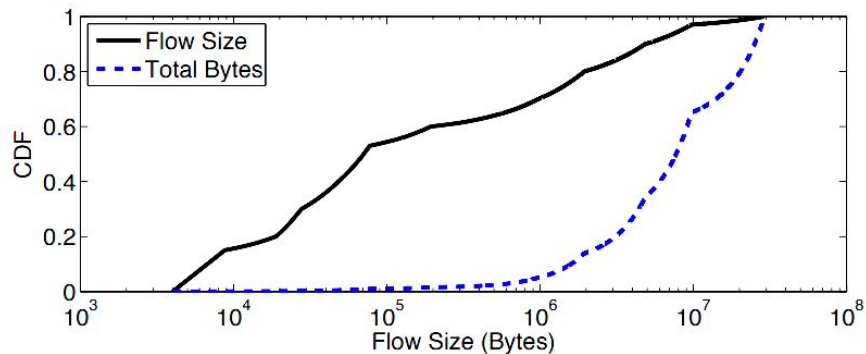# Rare glimpse of Google datacenter
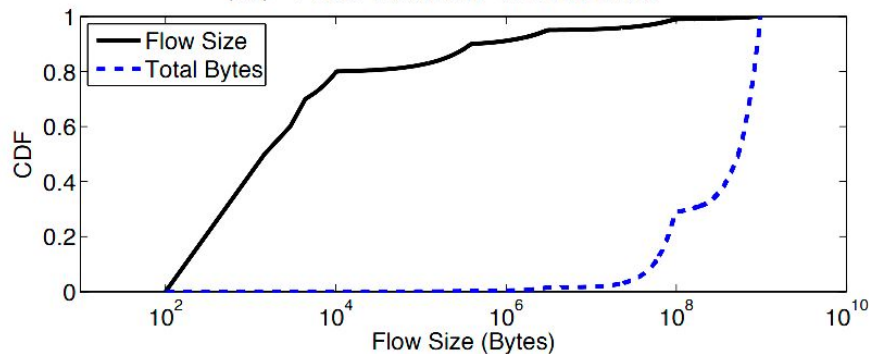
# Rare glimpse of Google datacenter

# Fear of the buffer

## Buffer per unit capacity (KB/Gbps)



Chart: KB per Gbps

| Br. 56538 | Trident+ | Trident II | Tomahawk | Tofino |
|-----------|----------|------------|----------|--------|
| ~85 | ~20 | ~10 | ~5 | ~3 |

# DC workloads and short flows



(a) Web search workload

(b) Data mining workload

# DC workloads and short flows



(a) Web search workload

(b) Data mining workload

Majority traffic volume is from long flows

Majority Flows are short

POWERTCP

# DC workloads and short flows
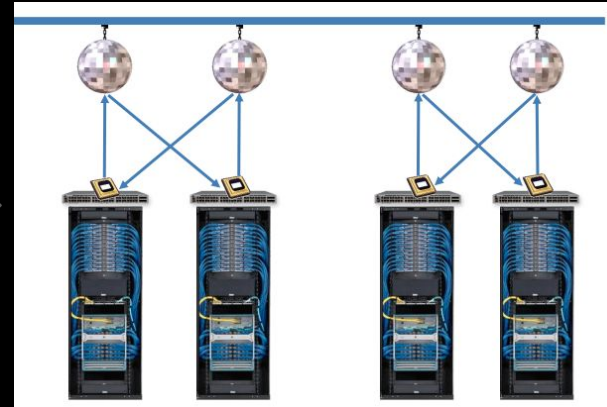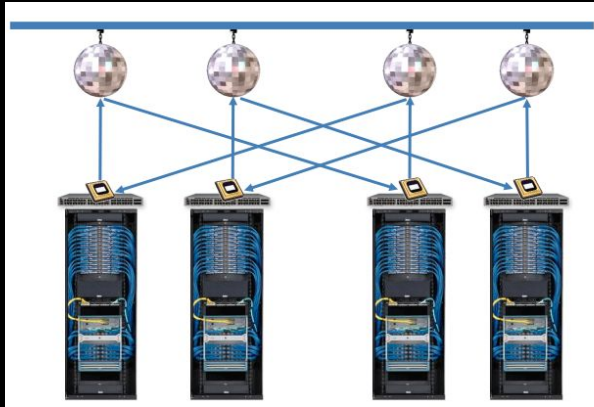

(b) Data mining workload

I have a phobia that throughput is always low

I have a constant fear that delay is always high

POWERTCP

# Emerging technologies and challenges

Not just queueing but quickly utilizing available bandwidth is important too

eg., Emerging Reconfigurable Datacenter Networks (RDCNs)

# Fine-grained congestion control is important for datacenter performance

# Timeline of congestion control in datacenters

- Reno, Cubic
- DCTCP, DCQCN
- Timely
- HPCC
- Swift

# Timeline of congestion control in datacenters

- **Voltage-based** (BDP + Queue Length)
  - ECN/Loss (*eg.*, DCTCP)
  - RTT based (*eg.*, Swift)
  - Inflight based (*eg.*, HPCC)
- **Current-based** (Total transmission rate)
  - RTT-gradient based (Eg., Timely)

PowerTCP

**Voltage-based**

Reaction to queue length or RTT

**Loss/ECN**
*eg., DCTCP*

**Voltage-based**

**Reaction to queue length or RTT**

POWERTCP

**Loss/ECN**
*eg., DCTCP*

**Delay**
*eg., Swift*

**Voltage-based**

**Reaction to queue length or RTT**

**Loss/ECN**
*eg., DCTCP*

**Delay**
*eg., Swift*

**Inflight**
*eg., HPCC*

**Voltage-based**

**Reaction to queue length or RTT**

POWERTCP

**Current-based**

Reaction to variations

RTT gradient
*eg., Timely*

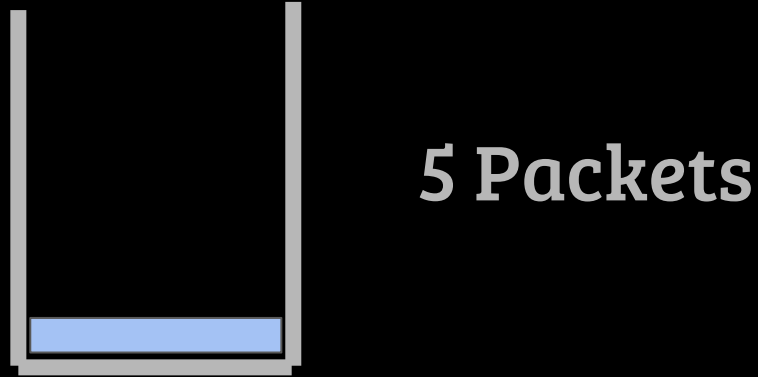Loss/ECN
*eg., DCTCP*

Delay
*eg., Swift*

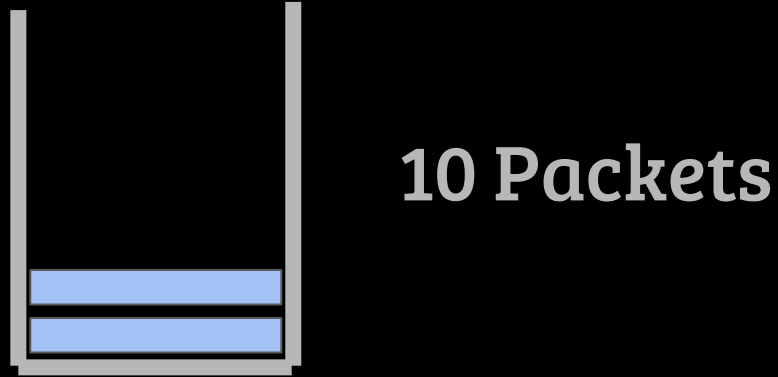Inflight
*eg., HPCC*

**Voltage-based**

**Reaction to queue length or RTT**

POWERTCP

29

# Problems of existing approaches

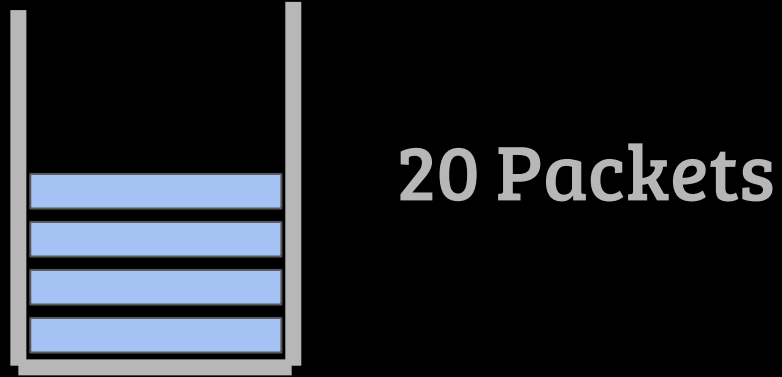**Fundamentally limited to a single dimension**

# Problems of existing approaches



5 Packets

# Problems of existing approaches



10 Packets

# Problems of existing approaches
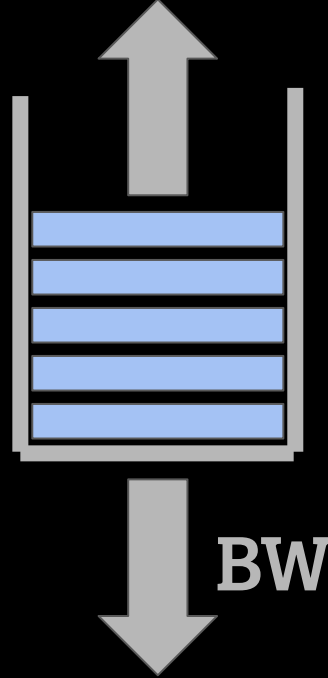


15 Packets

# Problems of existing approaches

20 Packets
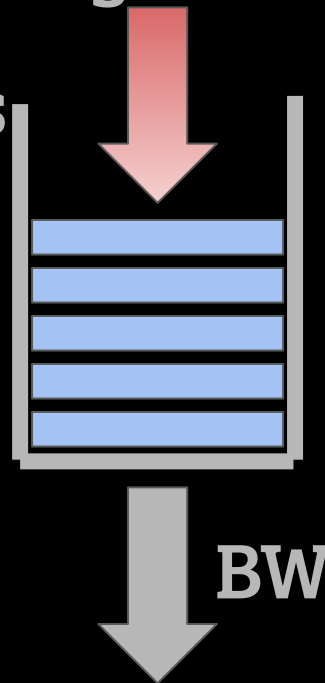
**Problems of existing approaches**

Increasing at 8  x BW

25 Packets

BW
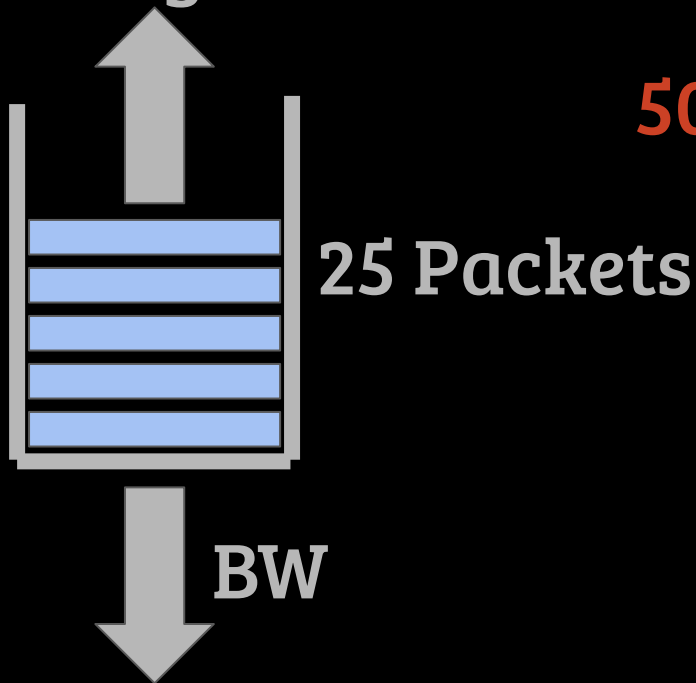
# Problems of existing approaches

## Increasing at 8x BW

25 Packets

BW

## Draining at max rate

25 Packets

BW

POWERTCP

# Problems of existing approaches

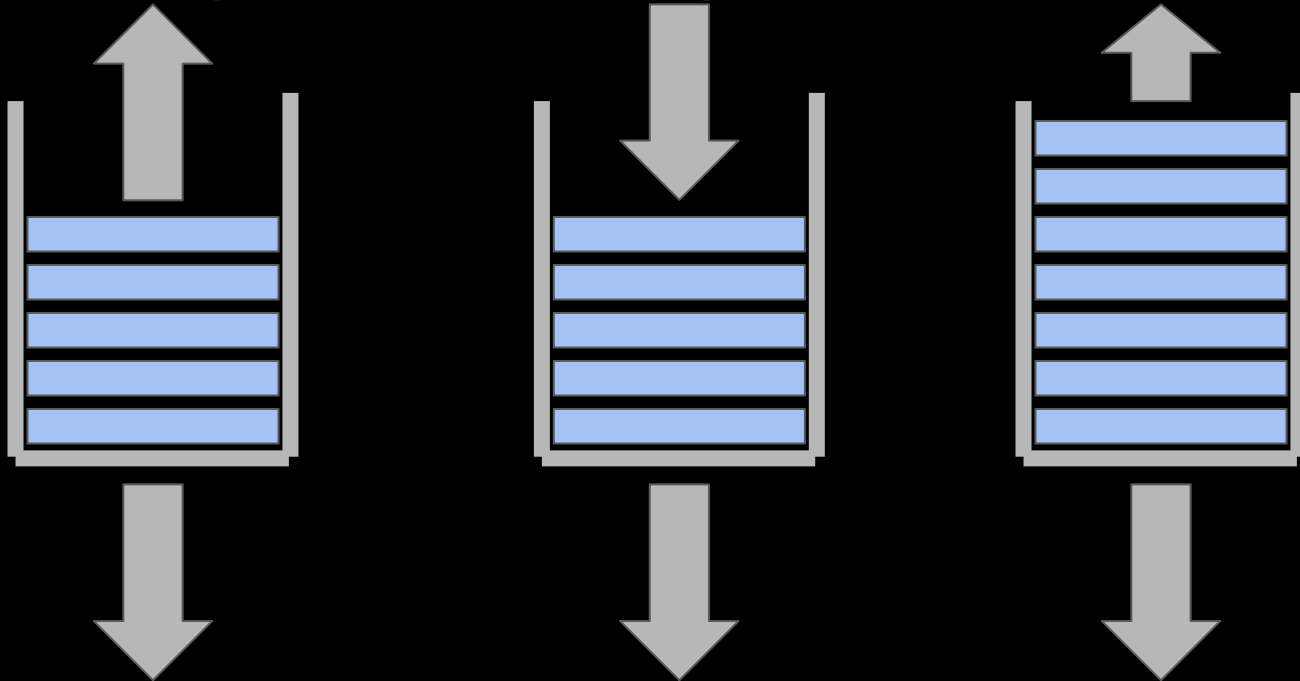## Increasing at 8x BW

## Increasing at 8x BW

**50 Packets**

25 Packets

BW

BW

# Problems of existing approaches

**Fundamentally limited to a single dimension**

# Summary of Our Analysis

- **Voltage-based**
    - Can in-principle achieve near-zero queue equilibrium
    - Slow reaction
- **Current-based**
    - Unstable with no equilibrium
    - Fast Reaction

Current-based

Reaction to variations

Timely

Better inflight control

DCTCP          Swift          HPCC

Voltage-based

Reaction to queue length or RTT

POWERTCP

**Current-based**

Reaction to variations

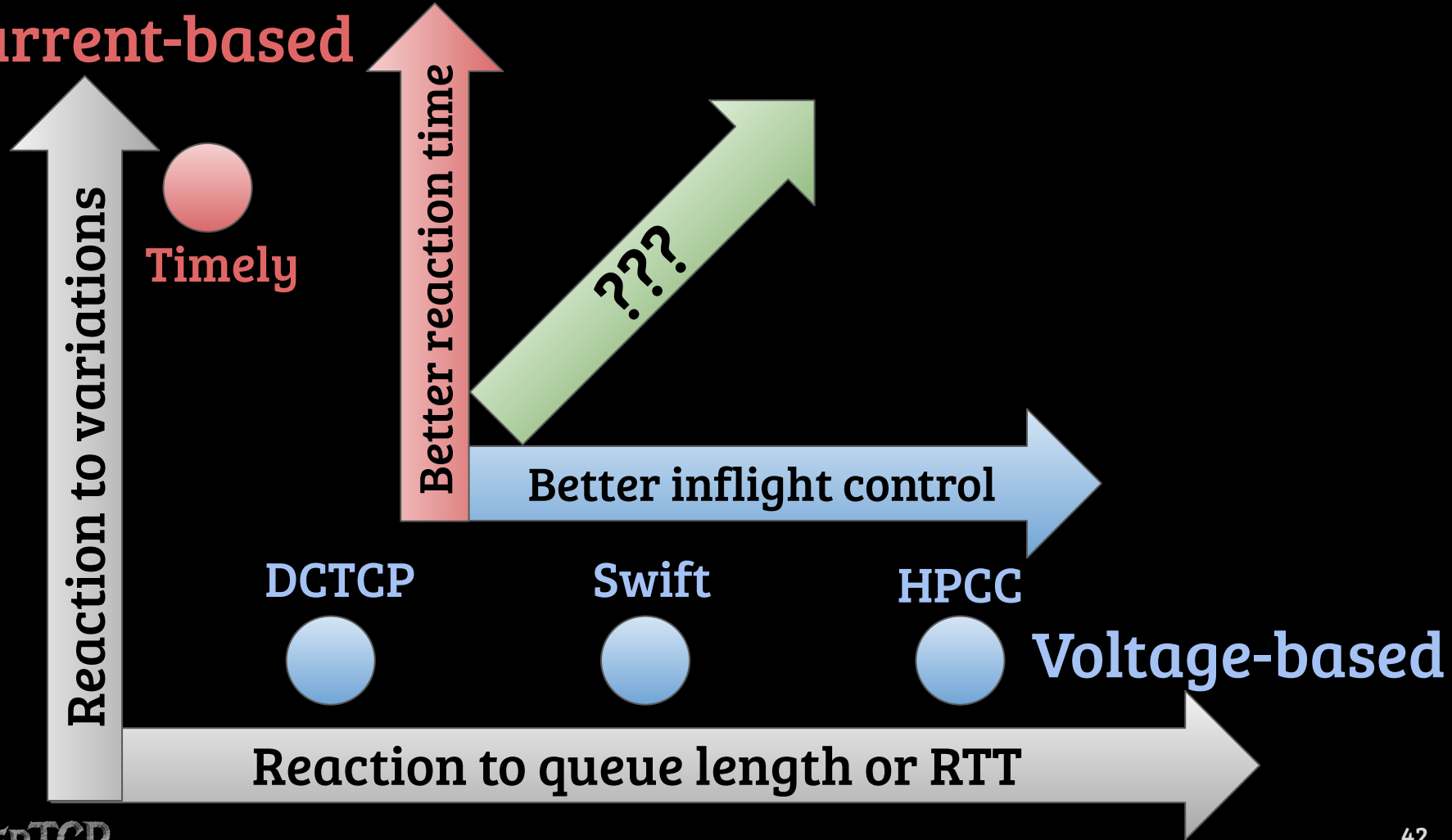Timely

Better reaction time

Better inflight control

DCTCP     Swift     HPCC

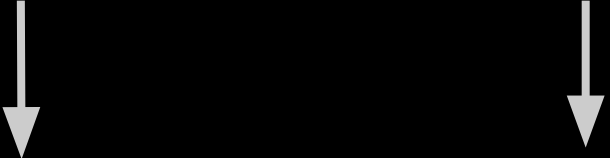**Voltage-based**

**Reaction to queue length or RTT**

POWERTCP

41

**Current-based**

Reaction to variations

Timely

Better reaction time

???

Better inflight control

DCTCP

Swift

HPCC

**Voltage-based**

Reaction to queue length or RTT

POWERTCP

# The notion of power

## Power = Voltage x Current

$$\underbrace{\Gamma}_{\text{Power}} = \underbrace{(q(t) + b \times \tau)}_{\text{Voltage}} \times \underbrace{(\dot{q}(t) + \mu(t))}_{\text{Current}}$$

Voltage → BDP+queue bytes

Current → Total rate

POWERTCP

# The notion of power

**Enqueue rate = queue-gradient + Dequeue rate**

$$\lambda(t - t^f) = \dot{q}(t) + \mu(t)$$

**Sending rate = Window per RTT**

$$\lambda(t) = \frac{w(t)}{\theta(t)}$$

**RTT = queueing delay + base RTT**

$$\theta(t - t^f) = \frac{q(t)}{b} + \tau$$

# The notion of power

$$b \times w(t - t^f) = \underbrace{(q(t) + b \times \tau)}_{\textbf{Voltage}} \times \underbrace{(\dot{q}(t) + \mu(t))}_{\textbf{Current}}$$

POWERTCP

# The notion of power

A function of both queue length and variations

# The notion of power

A function of both queue length and variations

- Detects increased queue lengths

# The notion of power

A function of both queue length and variations

- Detects increased queue lengths
- Detects congestion onset and intensity

# The notion of power

A function of both queue length and variations

- Detects increased queue lengths
- Detects congestion onset and intensity
- Detects rapid drop in queue lengths

Current-based

Reaction to variations

Timely

Better reaction time

Power-based CC

Better inflight control

DCTCP    Swift    HPCC

Voltage-based

Reaction to queue length or RTT

POWERTCP

# Current-based

Reaction to variations

Timely

Better reaction time

Power-based CC

**POWERTCP**

Better inflight control

DCTCP          Swift          HPCC          Voltage-based

**Reaction to queue length or RTT**

POWERTCP

# PowerTCP control law

$$w_i(t + \delta t) = \gamma \cdot \left( w_i(t) \cdot \frac{e}{f(t)} + \beta \right) + (1 - \gamma) \cdot w_i(t)$$

**New window size**

# PowerTCP control law

$$w_i(t + \delta t) = \gamma \cdot \left( w_i(t) \cdot \frac{e}{f(t)} + \beta \right) + (1 - \gamma) \cdot w_i(t)$$

**Old window size**

# PowerTCP control law

$$w_i(t + \delta t) = \gamma \cdot \left( w_i(t) \cdot \boxed{\frac{e}{f(t)}} + \beta \right) + (1 - \gamma) \cdot w_i(t)$$

**MIMD based on Power**
*(Multiplicative increase - multiplicative decrease)*

POWERTCP

# PowerTCP control law

$$w_i(t + \delta t) = \gamma \cdot \left( w_i(t) \cdot \frac{e}{f(t)} + \boxed{\beta} \right) + (1 - \gamma) \cdot w_i(t)$$

**Additive increase**

# PowerTCP control law

$$w_i(t + \delta t) = \boxed{\gamma} \cdot \left( w_i(t) \cdot \frac{e}{f(t)} + \beta \right) + \boxed{(1 - \gamma)} \cdot w_i(t)$$

**Exponential Weighted Moving Average (EWMA)**

# PowerTCP feedback

Power is measured via Inband Network Telemetry (INT)

- Queue lengths
- Timestamps
- Tx bytes
- Bandwidth

# PowerTCP without switch support

- Power can be measured via delay signal

POWERTCP

# PowerTCP without switch support

- Power can be measured via delay signal

$$\underbrace{\Gamma}_{\text{Power}} = b^2 \times \underbrace{\theta}_{\text{Voltage}} \times \underbrace{(\dot{\theta} + 1)}_{\text{Current}}$$

Voltage $\longrightarrow$ RTT

Current $\longrightarrow$ RTT gradient

# Evaluation

# Evaluation - Incast

# Evaluation - Incast

# Evaluation - Incast

# Evaluation - Incast

# Evaluation - Incast

# Evaluation - Fairness & Stability

# Evaluation - Fairness & Stability

# Evaluation - Fairness & Stability

# Evaluation - Fairness & Stability
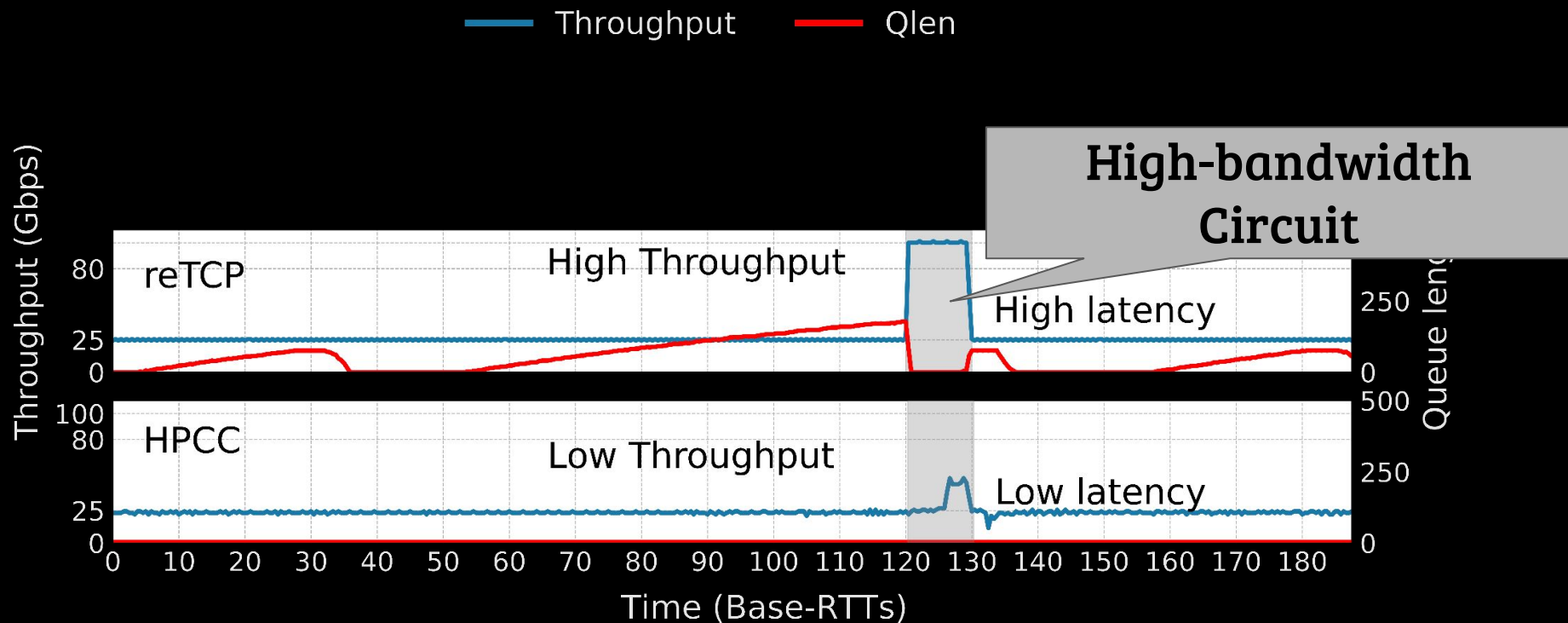
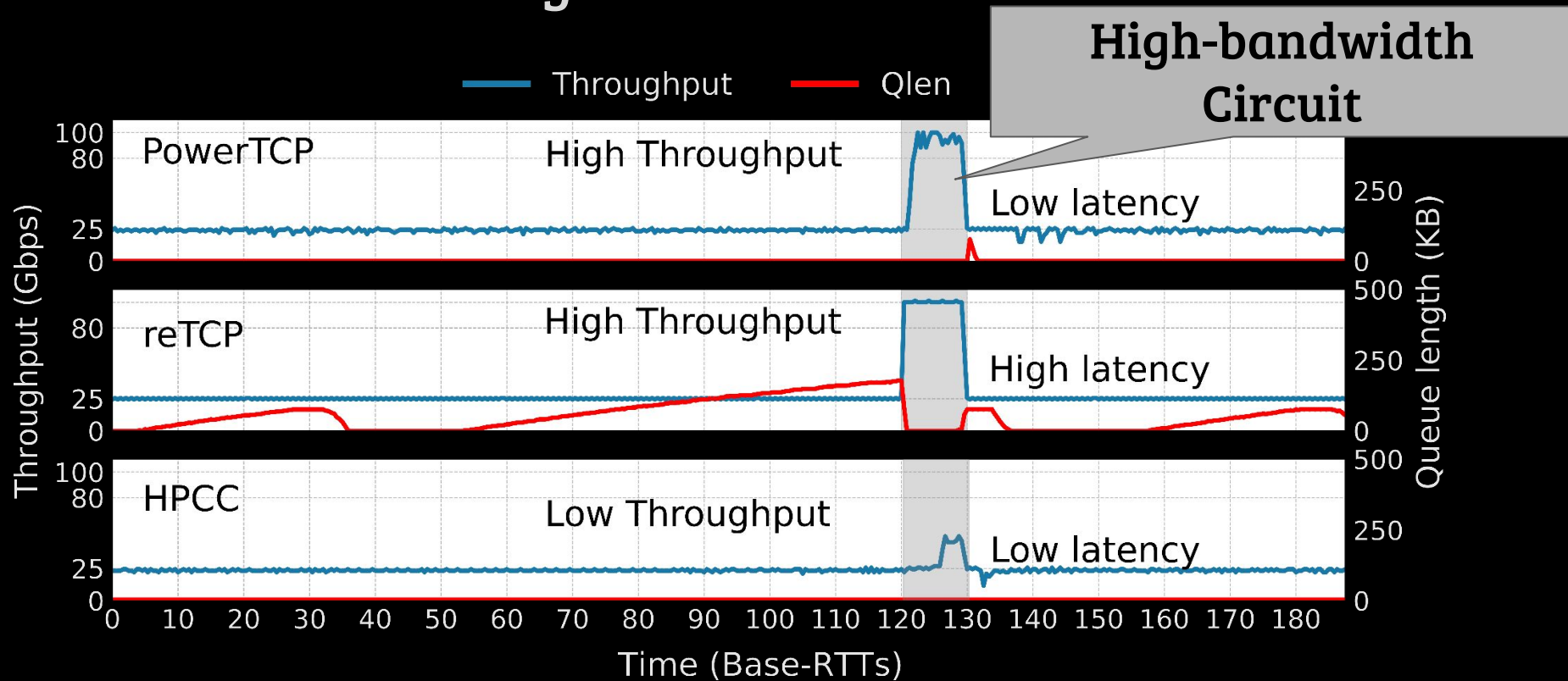# Evaluation - Workload

# Evaluation - Workload

# Evaluation - Reconfigurable Networks

# Evaluation - Reconfigurable Networks

# Evaluation - Reconfigurable Networks

# Conclusion

- Existing CC are fundamentally limited to a single dimension
- Power is an interesting and provably good measure for CC
- PowerTCP: a novel control law based on Power
- Improves FCTs for short flows and even for long flows

# Thank you